



Independent Component Analysis

Pierre Comon

► To cite this version:

Pierre Comon. Independent Component Analysis. J-L.Lacoume. Higher-Order Statistics, Elsevier, pp.29-38, 1992. hal-00346684

HAL Id: hal-00346684

<https://hal.science/hal-00346684>

Submitted on 12 Dec 2008

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Independent Component Analysis

Pierre COMON

THOMSON - SINTRA

Parc Sophia Antipolis, BP 138, 06561 Valbonne Cedex, France

Abstract

The Independent Component Analysis (ICA) of a random vector consists of searching the linear transformation that minimizes the statistical dependence between its components. In order to design a practical optimization criterion, the expansion of mutual information is being resorted to, as a function of cumulants. The concept of ICA may be seen as an extension of Principal Component Analysis, which only imposes independence up to the second order and consequently defines directions that are orthogonal. Applications of ICA include data compression, detection and localization of sources, or blind identification and deconvolution.

Contents

1. Introduction
 2. Statements related to statistical independence
 3. Optimization criteria
 4. A practical algorithm
 5. Conclusion
- References
 Appendices

1. Introduction

This paper attempts to provide a precise definition of ICA within an applicable mathematical framework. It is envisaged that this definition will provide a baseline for further development and application of the ICA concept.

In this discussion, the following linear statistical model is assumed:

$$(1) \quad y = Mx + v,$$

where x , y and v are random vectors with values in \mathbb{R}^N or \mathbb{C}^N and with zero mean and finite covariance, and M is

a regular square matrix. The problem set by ICA may be summarized as follows. Given realizations of y , it is desired to estimate both M and the corresponding realizations of x . To achieve this, we assume that the components of vector x are statistically independent. However, because of the presence of the noise v , it is in general impossible to recover exactly x , especially if v is non-gaussian. We shall however derive a process that delivers the best estimate z of x , according to an optimization criterion that we have named "contrast".

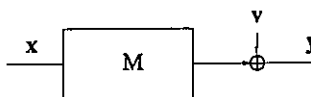


Figure 1

Let us dawdle a few lines on related applications. If x is a vector formed of N successive time samples of a white process, and if M is Toeplitz triangular, then model (1) represents nothing else but a deconvolution problem. The first column of M contains the successive samples of the impulse response of the corresponding FIR causal filter. Such blind identification and/or deconvolution problems have been addressed in different manners in [1] [12] [26] [27] [28] [32]. Note that if M is not triangular, the filter is allowed to be non-causal. Moreover, if M is not Toeplitz, the filter is allowed to be non-stationary.

In antenna processing, ICA may be utilized with the goal of localization of radiating sources with possibly perturbed or ill-calibrated arrays, as well as for detection and estimation of sources from unknown arrays without localization. The use of ICA has also some interesting features for jammer rejection, noise reduction, and blind equalization; it has been already experimented in radar for

instance [11]. Further ICA can be utilized in the identification of *multichannel* ARMA processes when the input is not observed and in particular for estimating the first coefficient of the model [8], which is generally assumed to be known [17]. On the other hand, ICA can be used as a data preprocessing tool before Bayesian detection and classification. In fact, by a change of coordinates, the density of multichannel data may be approximated by a product of marginal densities, allowing a density estimation with much shorter observations. Other related topics include high-order linear whitening and diagonalization of symmetric tensors.

The calculation of ICA was discussed in several recent papers [4] [8] [9] [13] [16] [18] [19] [23], where the problem was given various names. For instance the terminology "sources separation problem" has often been used. We shall not review existing works in detail here, but rather propose a steady definition. Our investigations have revealed that the problem of "Independent Component Analysis" was actually first proposed and so named by Herault and Jutten in 1986, refer to [19]. We shall retain this terminology. Furthermore, as demonstrated below, it has some links with Principal Component Analysis (PCA).

Suppose x is a non-degenerate gaussian random vector of dimension N with statistically independent components, and z a vector defined as $z = C x$, where C is a regular $N \times N$ matrix. Then if z has independent components, matrix C can be shown to be of the form:

$$(2) \quad C = \Lambda Q \Delta,$$

where both Λ and Δ are diagonal and Q orthogonal. This demonstrates that, if both x and z have a unit covariance matrix, then C may be any orthogonal matrix. We shall show later (theorem 23) that when x has at most one gaussian component, this indetermination reduces to a matrix of the form ΛP , where Λ is a diagonal matrix with entries of unit modulus, and P a permutation. This latter indetermination cannot be reduced further without additional assumptions. Note that PCA has exactly the same inherent indeterminations as ICA, so that we may assume the same arbitrary constraints in order to define ICA uniquely.

(3) DEFINITION

The ICA of a random vector y with finite covariance is a pair $\{F, \Delta\}$ of matrices such that:

(3a) the covariance of y decomposes into

$$V_y = F \Delta^2 F^*,$$

where Δ is diagonal real positive and F is full column rank;

(3b) Δ^2 is the covariance of a random vector z whose components are "the most independent possible", in the sense of the maximization of a given "contrast function".

(3c) The entries of Δ are sorted in decreasing order;

(3d) the columns of F are of unit norm;

(3e) the entry of largest modulus in each column of F is given a positive real part.

In this definition, superscript $(*)$ denotes transposition, and complex conjugation if the quantity is complex. As a consequence, ICA defines the so-called "source vector" z satisfying:

$$(4) \quad y = F z.$$

The requirement that F is square is not absolutely necessary, F could have more rows than columns. Nevertheless for the sake of clarity, this case is not discussed. As we shall see with theorem (23), ICA is unique as long as y has at most one gaussian component.

2. Statements related to independence

In this section, we shall first propose an appropriate contrast criterion. Then we state two results. It will be proved that ICA is uniquely defined if at most one component of x is gaussian, and then it will be shown why pairwise independence is a sufficient measure of statistical independence in our problem.

Most of the results presented in this paper hold true either for real or complex variables. However, some derivations would become much more complicated if derived in the complex case. Therefore only real variables will be considered for the sake of clarity. In the remaining, plain lowercase (resp. uppercase) letters denote in general scalar quantities (resp. tables with at least two indices, namely tensors), whereas boldface

lowercase letters denote column vectors with values in \mathbb{R}^N .

Let \mathbf{x} be a random variable with values in \mathbb{R}^N and denote $p_{\mathbf{x}}(\mathbf{u})$ its probability density function (pdf). Vector \mathbf{x} has mutually independent components if

$$(5) \quad p_{\mathbf{x}}(\mathbf{u}) = \prod_{i=1}^N p_{x_i}(u_i).$$

So a natural way of checking whether \mathbf{x} has independent components is to measure a distance between both sides of (5):

$$(6) \quad \delta(p_{\mathbf{x}}, \prod p_{x_i}).$$

In statistics, the large class of f -divergences is of key importance among the possible distance measures available [2]. In these measures the roles played by both densities are not always symmetric, so that we are not dealing with proper distances. For instance, the Kullback divergence is defined as:

$$(7) \quad \delta(p_{\mathbf{x}}, p_{\mathbf{z}}) = \int p_{\mathbf{x}}(\mathbf{u}) \log \frac{p_{\mathbf{x}}(\mathbf{u})}{p_{\mathbf{z}}(\mathbf{u})} d\mathbf{u}.$$

Recall that the Kullback divergence satisfies

$$(8) \quad \delta(p_{\mathbf{x}}, p_{\mathbf{z}}) \geq 0,$$

with equality if and only if $p_{\mathbf{x}}(\mathbf{u}) = p_{\mathbf{z}}(\mathbf{u})$ almost everywhere. This property results from the well-known inequality

$$\log w \leq w - 1, \text{ with equality iff } w = 1.$$

Now, if we look at the form of the Kullback divergence of (6), we obtain precisely the average mutual information of \mathbf{x} :

$$(9) \quad I(p_{\mathbf{x}}) = \int p_{\mathbf{x}}(\mathbf{u}) \log \frac{p_{\mathbf{x}}(\mathbf{u})}{\prod p_{x_i}(u_i)} d\mathbf{u}, \quad \mathbf{u} \in \mathbb{R}^N.$$

From (8), the mutual information cancels if and only if the variables x_i are mutually independent, and is strictly positive otherwise.

On the other hand define the differential entropy of \mathbf{x} as:

$$(10) \quad S(p_{\mathbf{x}}) = - \int p_{\mathbf{x}}(\mathbf{u}) \log p_{\mathbf{x}}(\mathbf{u}) d\mathbf{u}.$$

Remind that differential entropy is not the limit of Shannon's entropy defined for discrete variables; it is not invariant by change of coordinates as the entropy was. Yet, it is the usual practice to still call it entropy, in short. Entropy enjoys very privileged properties as emphasized in [29], and we shall show next that information (9) may also be written as a difference of entropies.

Now denote \mathbb{E}^N the space of random variables with values in \mathbb{R}^N , \mathbb{E}_2^N the Euclidian subspace of \mathbb{E}^N spanned by variables with finite covariance matrix, provided with the scalar product $\langle \mathbf{x}, \mathbf{y} \rangle = E\{\mathbf{x}^* \mathbf{y}\}$, and $\tilde{\mathbb{E}}_2^N$ the subset of \mathbb{E}_2^N of variables having a regular covariance. Lastly define \mathbb{E}_r^N as the euclidian subspace of \mathbb{E}_2^N constituted of variables with finite moments up to order r . Note that any random variable of \mathbb{E}^N , with finite moments up to order r and with a pdf not reduced to a point-like mass, can be reduced by projection to a variable belonging to $\tilde{\mathbb{E}}_r^N = \mathbb{E}_r^N \cap \tilde{\mathbb{E}}_2^N$ for some $N, N' \geq N \geq 1$.

Among the densities of $\tilde{\mathbb{E}}_2^N$ having a given covariance matrix V , the gaussian density is the one which has the largest entropy. If \mathbf{x} is zero-mean gaussian, its pdf will be referred to with the notation $\phi_{\mathbf{x}}(\mathbf{u})$, with:

$$(11) \quad \phi_{\mathbf{x}}(\mathbf{u}) = (2\pi)^{-N/2} |V|^{-1/2} \exp\{-\mathbf{u}^* V^{-1} \mathbf{u}\}/2.$$

Our proposition says that

$$(12) \quad S(\phi_{\mathbf{x}}) \geq S(p_{\mathbf{y}}),$$

with equality iff $\phi_{\mathbf{x}}(\mathbf{u}) = p_{\mathbf{y}}(\mathbf{u})$ almost everywhere. The entropy obtained in the case of equality is

$$(13) \quad S(\phi_{\mathbf{x}}) = \frac{1}{2} [N + N \log(2\pi) + \log \det V].$$

Other simple densities enjoy similar properties. For instance, the uniform density maximizes entropy over the class of densities with a bounded support. On the other hand, the exponential density maximizes entropy over the class of densities defined in the first quadrant, $(\mathbb{R}^+)^N$, and with a given mean. In the rest of the paper, there will be no restriction on the support of the densities, so that they will be defined on the entire space \mathbb{R}^N .

For densities in \mathbb{E}_2^N , one defines the negentropy as:

$$(14) \quad J(p_{\mathbf{x}}) = S(\phi_{\mathbf{x}}) - S(p_{\mathbf{x}}),$$

where $\phi_{\mathbf{x}}$ stands for the gaussian density with the same mean and variance as $p_{\mathbf{x}}$. Negentropy may be written in another manner, as a Kullback divergence:

$$(15) \quad J(p_{\mathbf{x}}) = \int p_{\mathbf{x}}(\mathbf{u}) \log \frac{p_{\mathbf{x}}(\mathbf{u})}{\phi_{\mathbf{x}}(\mathbf{u})} d\mathbf{u},$$

which shows, referring to (8), that

$$(16) \quad J(p_{\mathbf{x}}) \geq 0,$$

with equality iff $\phi_{\mathbf{x}} = p_{\mathbf{x}}$ almost everywhere.

From (9) and (15), the mutual information may be written as the difference of negentropies

$$(17) \quad I(p_x) = J(p_x) - \sum_{i=1}^N J(p_{x_i}) + \frac{1}{2} \ln \frac{\prod V_{ii}}{\det V},$$

where V denotes the variance of x . The proof is deferred to [7] for reasons of space. This relation gives a means of approximating the mutual information, provided we are able to approximate the negentropy about zero, which amounts to expanding the density p_x in the neighborhood of ϕ_x . This will be the starting point of section 3.

So far, it has been seen that the gaussian density plays a key role in the problem under consideration, since ICA is undetermined if p_x is gaussian. It is therefore natural to resort to negentropy which measures deviations from gaussianity. We have shown that both the gaussian feature and the mutual independence can be characterized with the help of entropy. Yet, these remarks justify only in part the use of (17) as an optimization criterion in our problem. In fact from (3), this criterion should meet the requirements given below.

(18) DEFINITION

A contrast is a mapping Ψ from \mathbb{E}^N to \mathbb{R} satisfying the 3 requirements:

- $\Psi(x)$ depends only on p_x , $\forall x \in \mathbb{E}^N$.
- Ψ is invariant by scale change, that is:

$$\Psi(\Lambda x) = \Psi(x), \forall \Lambda \text{ diagonal regular.}$$
- if x has independent components, then:

$$\Psi(A x) \leq \Psi(x), \forall A \text{ regular.}$$

Contrary to the criterion proposed in the subsequent theorem, it is easy to see that the information (17) is not scale invariant. Let z be a zero-mean random variable of $\tilde{\mathbb{E}}_2^N$, V its covariance matrix, and L a matrix such that $V = LL^*$ (it could be a Cholesky factorization, or any other square-root decomposition, based on SVD for instance). Then we define the standardized variable associated with z :

$$(19) \quad \tilde{z} = L^{-1} z.$$

Note that $\tilde{z}_i \neq \tilde{z}_i$. In fact, \tilde{z}_i is merely the variable z_i normalized by its variance. In the following, we shall only talk about \tilde{z}_i , the i th component of the standardized variable \tilde{z} . It is easy to see, by the way, that entropies of z and \tilde{z} are related by [7]

$$(20) \quad S(p_z) = S(p_{\tilde{z}}) - \frac{1}{2} \log \det V.$$

Now we are in a position to define a contrast criterion.

(21) THEOREM

The following mapping is a contrast over $\tilde{\mathbb{E}}_2^N$:

$$\Psi(p_z) = -I(p_{\tilde{z}}).$$

See [7] for a proof. Note that from (17) we have

$$\Psi(p_z) = \sum_{i=1}^N J(p_{x_i}) - J(p_x).$$

Criterion (21) is admissible for ICA computation. This theoretical criterion, involving a generally unknown density, will be made usable by approximations in section 3. Regarding computational loads, the calculation of ICA may still be too heavy even after approximations, and we now turn to a theorem that theoretically explains why the practical algorithm designed in section 4, that proceeds pairwise, indeed works.

(22) LEMMA

Let x and z be two random vectors such that $z = B x$, where B is regular. Suppose additionally that x has independent components, and that z has pairwise independent components. If B has two non-zero entries in the same column j , then x_j is either gaussian or constant.

(23) THEOREM

Let x be a vector with independent components of which at most one is gaussian, and whose densities are not reduced to a point-like mass. Let C be an orthogonal $N \times N$ matrix and z the vector $z = C x$. Then the 3 following properties are equivalent:

- (i) The components z_i are pairwise independent
- (ii) the components z_i are mutually independent
- (iii) $C = \Lambda P$, Λ diagonal, P permutation.

See appendix for a proof of lemma and theorem.

3. Optimization criteria

Suppose that we observe \tilde{y} and that we look for an orthogonal matrix Q maximizing the contrast:

$$(24) \quad \Psi(p_z) = -I(p_{Q\tilde{y}}),$$

where

$$\tilde{z} = Q \tilde{y}.$$

In practice, the densities $p_{\tilde{z}}$ and $p_{\tilde{y}}$ are not known, so that

the criterion (24) cannot be directly utilized. The aim of this section is to express the contrast (21) as a function of the standardized cumulants (of order 3 and 4), which are quantities more easily accessible. The expression of entropy and negentropy in the scalar case will be first briefly derived. We start with the Edgeworth expansion of type A of a density. A central limit theorem says that if z is a sum of P independent random variables with finite cumulants, then the i th order cumulant of z is of order:

$$(25) \quad \kappa_i \sim P^{\frac{2-i}{2}}.$$

This theorem can be traced back to 1928 and is attributed to Cramér [33]. Referring to [20, p176, formula 6.49], the expansion of the pdf of z up to order 4 about its best gaussian approximate (here with zero-mean and unit variance) is given by

$$(26) \quad \frac{p_z(u)}{\phi_z(u)} = 1 + \frac{1}{3!} \kappa_3 h_3(u) + \frac{1}{4!} \kappa_4 h_4(u) + \frac{10}{6!} \kappa_3^2 h_6(u) + \frac{1}{5!} \kappa_5 h_5(u) + \frac{35}{7!} \kappa_3 \kappa_4 h_7(u) + \frac{280}{9!} \kappa_3^3 h_9(u) + \frac{1}{6!} \kappa_6 h_6(u) + \frac{56}{8!} \kappa_3 \kappa_5 h_8(u) + \frac{35}{8!} \kappa_4^2 h_8(u) + \frac{2100}{10!} \kappa_3^2 \kappa_4 h_{10}(u) + \frac{15400}{12!} \kappa_3^4 h_{12}(u) + o(P^{-2}).$$

In this expression, κ_i denotes the cumulant of order i of the standardized scalar variable considered (this is the notation of Kendall and Stuart, not assumed subsequently in the multichannel case), and $h_i(u)$ is the Hermite polynomial of degree i . The advantage of Edgeworth expansion over Gram-Charlier's lies in the ordering of terms according to their decreasing significance as a function of $P^{-1/2}$. See [20] [21] for general remarks on pdf expansions.

Now let us turn to the expansion of the negentropy defined in (15). We start with

$$(27) \quad (1+v) \log(1+v) = v + v^2/2 - v^3/6 + v^4/12 + o(v^4),$$

and with the properties satisfied by the $h_i(u)$'s [7]:

$$(28) \quad \int \phi(u) h_i(u) du = 0,$$

$$(29) \quad \int \phi(u) h_p(u) h_q(u) du = \delta_{pq} p!,$$

$$(30) \quad \int \phi(u) h_3^2(u) h_4(u) du = 3!^3,$$

$$(31) \quad \int \phi(u) h_3^2(u) h_6(u) du = 6!,$$

$$(32) \quad \int \phi(u) h_3^4(u) du = 93 \cdot 3!^2.$$

THEOREM

Using these relations together with (26) and (27), one can prove that for a standardized scalar variable z :

$$(33) \quad J(p_z) = \frac{1}{12} \kappa_3^2 + \frac{1}{48} \kappa_4^2 + \frac{7}{48} \kappa_3^4 - \frac{1}{8} \kappa_3^2 \kappa_4 + o(P^{-2}).$$

See [7] for a proof. Next, from (17), the calculus of the mutual information of a standardized variable \tilde{z} needs not only the marginal negentropy of each component \tilde{z}_i but also the joint negentropy of \tilde{z} . It turns out that the calculation of $J(p_{\tilde{z}})$ is much more complicated than for the scalar case given above, though it goes along the same lines. Assume the Einstein summing convention, where the presence of the same index on top and bottom means a summation over this index, e.g.

$$K^{ijk} K_{ij} K_k = \sum_{ijk} K_{ijk} K_{ij} K_k.$$

This notation has some important mathematical meaning when used in the relevant context [5] [22] [6]. Then after very cumbersome and tedious calculations, one can show that¹

$$(34) \quad J(p_{\tilde{z}}) = \frac{1}{12} K^{ijk} K_{ijk} + \frac{1}{48} K^{ijkl} K_{ijkl} + \frac{3}{48} K^{ijk} K_{ijn} K_{kqr} K^{qrm} + \frac{4}{48} K^{ijk} K_{imn} K_{jr}^m K_k^{nr} - \frac{1}{8} K^{ijk} K_{ilm} K_{jk}^{lm} + o(P^{-2}),$$

where $K_{ij\dots q}$ denote the cumulants $\text{Cum}\{\tilde{z}_i, \tilde{z}_j, \dots, \tilde{z}_q\}$. However, it is more clever to notice that

$$(35) \quad J(p_{\tilde{z}}) = J(p_{\tilde{y}})$$

since entropy is invariant by orthogonal change of coordinates. Using (33) and (35), the mutual information of \tilde{z} takes the form, up to $O(P^{-2})$ terms:

$$(36) \quad I(p_{\tilde{z}}) = J(p_{\tilde{y}}) - \frac{1}{48} \sum_{i=1}^N \{4 K_{iii}^2 + K_{iii}^2 + 7 K_{iii}^4 - 6 K_{iii}^2 K_{iii}\}.$$

Yet, cumulants satisfy a multilinearity property [3], which allows them to be called tensors [24] [25]. Denote K the family of cumulants of \tilde{z} and Γ that of \tilde{y} . Then, by resorting to Einstein notation for the sake of conciseness, this property can be written at orders 3 and 4 as:

¹: the exact form of this expansion was derived with the help of J.F.Podevin, a student visiting the author.

$$(37) \quad K_{ijk} = Q_i^p Q_j^q Q_k^r \Gamma_{pqr},$$

$$(38) \quad K_{ijkl} = Q_i^p Q_j^q Q_k^r Q_l^s \Gamma_{pqrs}.$$

On the other hand, $J(p_y)$ does not depend on Q so that the criterion (24) can be reduced up to order P^2 to the maximization of a functional ψ :

$$(39) \quad \psi(Q) = \sum_{i=1}^N 4 K_{iii}^2 + K_{iii}^2 + 7 K_{iii}^4 - 6 K_{iii}^2 K_{iiii}$$

with respect to Q , keeping in mind that the tensors K depend on Q through relations (37) and (38). The function $\psi(Q)$ is actually a complicated rational function in $N(N-1)/2$ variables. The goal of the remainder of the paper is to avoid exhaustive search and save computational time in the optimization problem.

Simpler criteria

(41) LEMMA

Denote rQ the matrix obtained by raising each entry of an orthogonal matrix Q to the power r . Then we have:

$$\| {}^2Q \, u \| \leq \| u \|.$$

(42) THEOREM

The functional

$$\psi(Q) = \sum_{i=1}^N K_{ii \dots i}^2$$

where $K_{ii \dots i}$ are marginal standardized cumulants of order r , is a contrast for any $r > 2$.

The proofs are given in appendix. See also [10] or [7] for more details. These contrast functions are generally less discriminating than (24). In fact, if two components have a zero cumulant of order r , the contrast (42) fails to separate them (this is the same behavior as for gaussian components). However, as in theorem (23), at most one source component is allowed to have a null cumulant.

Now, it is easy to show [10] [7] that the quantity

$$(43) \quad \Omega_r = \sum_{i_1 \dots i_r} K_{i_1 i_2 \dots i_r}^2$$

is invariant under linear and regular transformations. This result gives an interpretation to the significance of contrast functions such as (42). Indeed, the maximization of $\psi(Q)$ is equivalent to the minimization of $\Omega_r - \psi(Q)$, which is eventually the same as to minimize the sum of the squares of all cross-cumulants of order r , and these cumulants are precisely the measure of statistical dependence at order r . The same interpretation can be

given for the contrast (39) since the expression

$$(44) \quad \Omega_{3,4} = \sum_{ijklmnqr} 4 K_{ijk}^2 + K_{ijkl}^2 + 3 K_{ijk} K_{ijn} K_{kqr} K_{qrm} \\ + 4 K_{ijk} K_{imn} K_{jmr} K_{knr} - 6 K_{ijk} K_{ilm} K_{jklm}$$

is invariant under linear and regular transformations [7].

The interest in maximizing the contrast (42) rather than minimizing cross cumulants lies essentially in the fact that only N cumulants are involved instead of $O(N^2)$. Thus, this spares a lot of computations when estimating the cumulants from the data. A first analysis of complexity was given in [6] and [10].

Link with blind identification and deconvolution

The criteria (39) and (42) may be connected with other criteria recently proposed in the literature for the blind deconvolution problem. For instance, the criterion proposed in [28] may be seen to be equivalent to maximize $\sum K_{iii}^2$. In [12], one of the optimization criteria proposed amounts to minimizing $\sum S(p_{z_i})$, which is consistent with (14), (17) and (24). In [1], the family of criteria proposed contains (24).

On the other hand, identification techniques presented in [31] [26] [27] [17] [30] solve a system of equations obtained by cumulant matching. Though they work quite well in general, these approaches may seem arbitrary in their selection of particular equations rather than others. Moreover, their robustness are questioned in presence of measurement noise, especially non-gaussian, as well as for short data records. In [32] a matching in the Least Squares (LS) sense is proposed; in [8] the use of much more equations than unknowns improves on robustness for short data records. A more general approach is developed in [15] where a weighted LS matching is suggested. Our equation (39) gives some justifications to the process of selecting particular cumulants, by showing their dominance faced to the others. In this context, some simplifications would occur when developing (39) as a function of Q since the components y_i are generally assumed identically distributed (this is not assumed in our derivations). We insist however that the fourth order cumulant cannot be isolated from the third order ones, except if they all vanish.

4. A practical algorithm

Pairwise processing

As suggested by theorem (23), in order to maximize (39) it is necessary and sufficient to consider only pairwise cumulants of \tilde{z} . The proof of (23) was valid for any contrast but only in the case where the observation y was effectively stemming linearly from a random variable x with independent components (model (1) in the noiseless case). It turns out that it is true for any \tilde{y} and any $\tilde{z} = Q \tilde{y}$, and for any contrast of polynomial form in marginal cumulants of \tilde{z} as it is the case in (39) or (42). The proof resorts to differentiation tools borrowed from classical analysis, and not to statistical considerations [7]. These statements justify a pairwise processing.

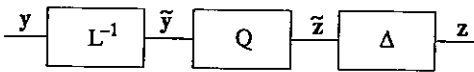


Figure 2: processing scheme in three main steps. Two additional transforms, P and Λ , must be applied afterwards in order to meet all the requirements of definition (3).

Given a set of data, $Y = \{y(t), 1 \leq t \leq T\}$, the proposed algorithm processes each pair in turn, similarly to the Jacobi algorithm in the diagonalization of symmetric real matrices:

(45) ALGORITHM

- 1) Compute a matrix L and the corresponding standardized data, $Z = L^{-1} Y$. L may be based on the QR factorization of Y or on its PCA.
- 2) Initialize $F = L$.
- 3) Begin a loop on the sweeps: $k=1, 2, \dots$
- 4) Sweep the $N(N-1)/2$ pairs (i, j) , according to a fixed ordering. For each pair:
 - a. estimate the required cumulants of $(Z_{i\cdot}, Z_{j\cdot})$, by resorting to κ -statistics for instance [20] [10].
 - b. find the angle θ maximizing $\psi(Q^{(i,j)})$, where $Q^{(i,j)}$ is the Givens rotation of angle θ , $\theta \in]-\pi/4, \pi/4[$.
 - c. accumulate $F := F Q^{(i,j)*}$.
 - d. update $Z := Q^{(i,j)} Z$.
- 5) End of the loop on k if $k=k_{\max}$ or if all estimated angles are very small.
- 6) Compute the norm of the columns of F : $\Delta_{ii} = \|F_{:,i}\|$.
- 7) Sort the entries of Δ in decreasing order:
 $\Delta := P \Delta P^*$ and $F := F P^*$.
- 8) Normalize F by the transform $F := F \Delta^{-1}$.

- 9) Fix the phase (sign) of each column of F according to (3e). This yields $F := F \Lambda$.

As shown in [10], the step 4)b can be carried out in various manners, and can become very simple when a contrast of type (42) is used. Real-time implementations are also possible [9]. See [10] and [7] for more details.

Robustness in presence of non-gaussian noise

In this section, the behavior of ICA in presence of non-gaussian noise is investigated by means of simulations. We need first to define a distance between matrices modulo a multiplicative factor of the form ΛP , where Λ is diagonal regular, and P is a permutation. Let A and \hat{A} be two regular matrices, and define the matrices with unit-norm columns

$$\underline{A} = A \Delta^{-1}, \hat{\underline{A}} = \hat{A} \hat{\Delta}^{-1}, \text{ with } \Delta_{kk} = \|A_{:,k}\|, \hat{\Delta}_{kk} = \|\hat{A}_{:,k}\|.$$

The gap $\varepsilon(A, \hat{A})$ is built from the matrix $D = \underline{A}^{-1} \hat{\underline{A}}$ as:

$$(46) \quad \varepsilon(A, \hat{A}) = \sum_i \left| \sum_j |D_{ij}| - 1 \right|^2 + \sum_j \left| \sum_i |D_{ij}| - 1 \right|^2 + \sum_i \left| \sum_j |D_{ij}|^2 - 1 \right| + \sum_j \left| \sum_i |D_{ij}|^2 - 1 \right|.$$

It can be shown [7] that this measure of distance is indeed invariant by postmultiplication by a matrix of the form ΛP : $\varepsilon(A \Lambda P, \hat{A}) = \varepsilon(A, \hat{A} \Lambda P) = \varepsilon(A, \hat{A})$.

Consider the observations in dimension $N=2$:

$$(47) \quad y(t) = M x(t) + \eta w(t),$$

where $1 \leq t \leq T$, x and w are zero-mean and standardized, x is formed of two independent random variables of kurtosis -1.2 and -1.5 respectively, w_i are uniformly distributed in $[-\sqrt{3}, \sqrt{3}]$ (and have thus a kurtosis of -1.2), η is a positive real parameter, and

$$M = \begin{pmatrix} 1 & 3 \\ -3 & 1 \end{pmatrix}.$$

Then when η increases, the signal to noise ratio decreases. We give² in figure 3 the behavior of algorithm (45) when the contrast (42) is used at order $r=4$ (which also coincides with (39) since the densities are symmetrically distributed in this simulation).

²: simulations presented in figures 3 and 4 have been performed with the help of D.Cren, a student who visited the author in 1990.

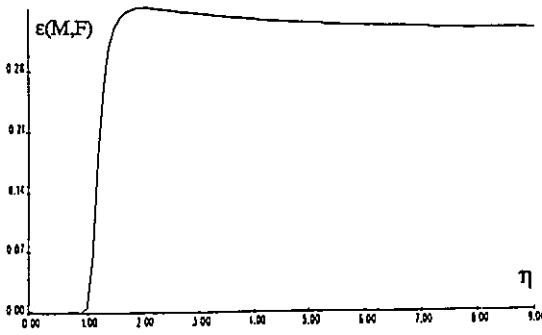


Figure 3: gap $\varepsilon(M,F)$ for 90 independent realizations of Y of size $T=2000$, each for 90 different values of η (by steps of 0.1).

The most surprising in these simulations is the excellent behavior of ICA even in the close neighborhood of $\eta = 1$. For $\eta > 1$, the ICA algorithm considers the vector Mx as a noise, and the vector ηw as a source vector. Because w has independent components, it could be checked that $\varepsilon(I,F)$ tends to zero as η increases to infinity, showing that matrix F is approaching AP .

Another interesting experiment is to look at the influence of the choice of the sweeping strategy on the convergence of the algorithm. For this purpose, consider now eleven independent sources. Here, we assume a fixed cyclic-by-rows description of the pairs, but the sources are shuffled. In ordering 1, the source kurtosis are

$$(1 -1 1 -1 1 -1 1 -1 1 -1 0)$$

whereas in ordering 2 they are

$$(1 1 1 1 1 -1 -1 -1 -1 -1 0).$$

Note the presence in this simulation of a null cumulant, and of cumulants of opposite signs. The mixing matrix is defined by $M_{ij} = 1 + \delta_{ij}$, and the additive noise has null kurtosis.

This simulation was performed directly from cumulants, so that the performances obtained are those that would be obtained for $T = \infty$ with real-world signals. The contrast utilized is (42) with $r=4$. It is plotted in figure 4(a) for ordering 1, and in figure 4(c) for ordering 2, as a function of the number of Givens rotations computed. For convenience, the gap between the original matrix, M , and the estimated matrix, $F = LQ^*\Delta^{-1}$ is also plotted in figures 4(b) and 4(d).

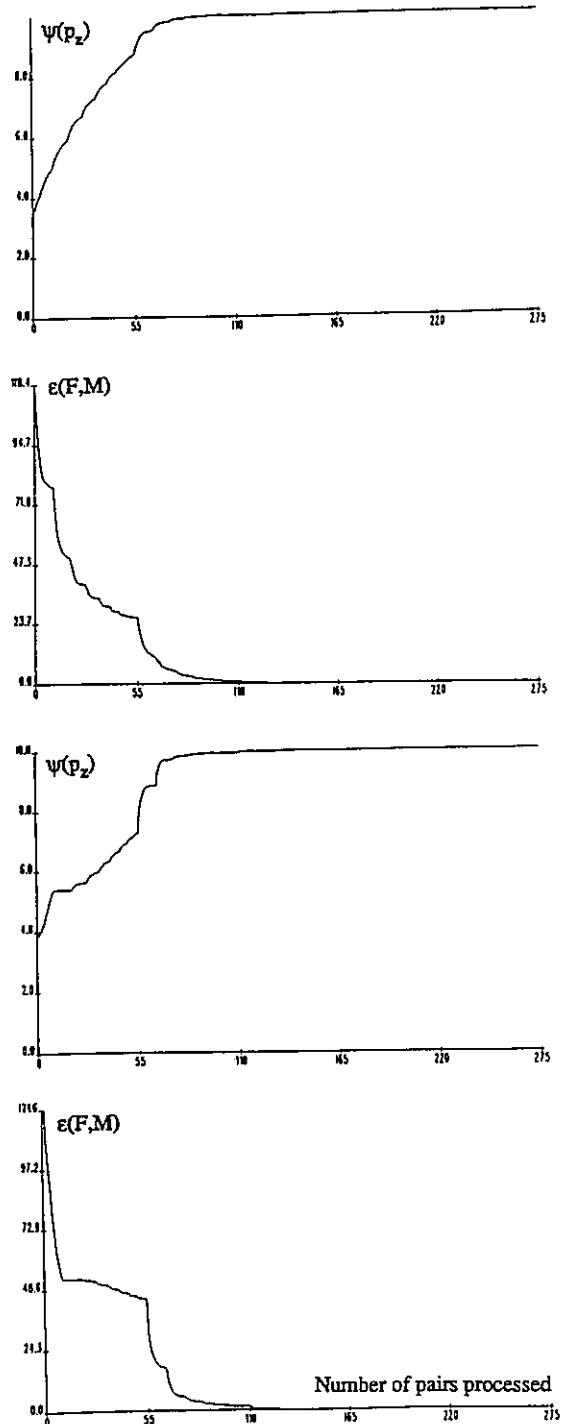


Figure 4: influence of the ordering of the sources on the convergence speed. The contrast ψ and the gap ε are plotted in (a) and (b) respectively for ordering 1, and in (c) and (d) for ordering 2.

In real-world experiments, the gap cannot be accessed so that the stopping criterion can only be based either on the contrast variation, or on the significance of the plane rotations.

Other simulations results related to ICA are reported in [10] and [7].

5. Conclusion

The definition of ICA given in this framework depends on a contrast function that serves as an optimization criterion. One of the contrasts proposed is built from the mutual information of standardized observations.

ICA can be computed by maximizing a combination of third and fourth order cumulants, obtained from the Edgeworth expansion of the mutual information. For the purposes of complexity reduction, an algorithm is proposed that minimizes a simplified criterion, by rooting a sequence of polynomials of degree four.

References

- [1] A.BENVENISTE, M.GOURSAT and G.RUGET, "Robust Identification of a NonMinimum Phase System", *IEEE Trans. Automatic Control*, vol.25, n°3, 1980, 385-399.
- [2] M.BASSEVILLE, "Distance Measures for Signal Processing and Pattern Recognition", *Signal Processing*, vol.18, n°4, dec 1989, 349-369.
- [3] D.R.BRILLINGER, *Time Series Data Analysis and Theory*, Holden Day, 1981.
- [4] J.F.CARDOSO, "Super-Symmetric Decomposition of the Fourth-Order Cumulant Tensor, Blind Identification of more Sources than Sensors", *Proc. ICASSP 91*, may 14-17, 1991.
- [5] J.F.CARDOSO, "High-Order Narrow-Band Array Processing", *Workshop High - Order Statistics*, Chamrousse, France, july 10-12, 1991.
- [6] J.F.CARDOSO and P.COMON, "Tensor-Based Independent Component Analysis", *Conference EUSIPCO*, sept 18-21, 1990, Barcelona, 673-676.
- [7] P.COMON, "Independent Component Analysis", submitted to *Signal Processing*, special issue on High-Order Statistics to appear in spring 1992.
- [8] P.COMON, "Separation of Sources using High-Order Cumulants", *SPIE Conference on Advanced Algorithms and Architectures for Signal Processing*, vol. Real-time signal processing XII, San Diego, aug 8-10, 1989, 170-181.
- [9] P.COMON, "Process and Device for Real-time Signals Separation", *Patent* registered for Thomson-Sintra, n°9000436, dec 1989.
- [10] P.COMON, "Independent Component Analysis and Blind Identification", *Traitement du Signal*, vol.7, n°5, dec 1990, 435-450.
- [11] G.DESODT and D.MULLER, "Complex Independent Component Analysis applied to the Separation of Radar Signals", *Proc. EUSIPCO Conference*, Barcelona, Torres-Masgrau-Lagunas editors, Elsevier Science Publ, 1990.
- [12] D.L.DONOHU, "On Minimum Entropy Deconvolution", *Proc of the 2nd Appl. Time Series Symp*, Tulsa, 1980, reprinted in *Applied Time Series Analysis II*, Academic Press, 1981, 565-608.
- [13] P.DUVAUT, "Principles of Sources Separation Methods based on Higher-Order Statistics", *Traitement du Signal*, vol.7, n°5, dec 1990, 407-418.
- [14] G.DARMOIS, "Analyse Generale des Liaisons Stochastiques", *Rev Inst Internat Stat*, vol.21, 1953, 2-8.
- [15] B.FRIEDLANDER and B.PORAT, "Asymptotically Optimal Estimation of MA and ARMA Parameters of Non-Gaussian Processes from High-Order Moments", *IEEE Trans Auto Control*, vol.35, jan 1990, 27-35.
- [16] M.GAETAT and J.L.LACOUME, "Source Separation without A Priori Knowledge: the Maximum Likelihood Solution", *Proc. EUSIPCO Conference*, Barcelona, Torres-Masgrau-Lagunas editors, Elsevier Science Publ, 1990.
- [17] G.GIANNAKIS, Y.INOUE and J.M.MENDEL, "Cumulant-based Identification of Multichannel Moving Average Models", *IEEE Automatic Control*, vol.34, july 1989, 783-787.
- [18] Y.INOUE and T.MATSUI, "Cumulant Based Parameter Estimation of Linear Systems", *Proc Workshop Higher-Order Spectral Analysis*, june 1989, Vail, Colorado, 180-185.
- [19] C.JUTTEN and J.HERAULT, "Blind Separation of Sources, Part I", to appear in *Signal Processing*, vol.24, n°1, sept 1991.
- [20] M.KENDALL and A.STUART, *The Advanced Theory of Statistics*, vol.1, 1977.
- [21] S.KOTZ and N.L.JOHNSON, *Encyclopedia of Statistical Sciences*, Wiley, 1982.
- [22] J.L.LACOUME, "Tensor-based Approach of Random Variables", talk given at ENST Paris, feb 7, 1991.
- [23] J.L.LACOUME and P.RUIZ, "Extraction of Independent Components from Correlated Inputs, a Solution based on Cumulants", *Proc Workshop Higher-Order Spectral Analysis*, june 1989, Vail, Colorado, 146-151.
- [24] P.McCULLAGH, "Tensor Notation and Cumulants", *Biometrika*, vol.71, 1984, 461-476.
- [25] P.McCULLAGH, *Tensor Methods in Statistics*, Chapman and Hall, 1987.
- [26] C.L.NIKIAS, "ARMA Bispectrum Approach to NonMinimum Phase System Identification", *IEEE Trans ASSP*, vol.36, april 1988, 513-524.
- [27] C.L.NIKIAS and M.R.RAGHUVEER, "Bispectrum Estimation: A Digital Signal Processing Framework",

Proceedings of the IEEE, vol.75, n°7, July 1987, 869-891.

[28] O.SHALVI and E.WEINSTEIN, "New Criteria for Blind Deconvolution of NonMinimum Phase Systems", *IEEE Trans. Information Theory*, vol.36, n°2, march 1990, 312-321.

[29] J.E.SHORE and R.W.JOHNSON, "Axiomatic Derivation of the Principle of Maximum Entropy", *IEEE Trans. Inf. Theory*, vol.26, jan 1980, 26-37.

[30] A.SWAMI and J.M.MENDEL, "ARMA Parameter Estimation using Only Output Cumulants", *IEEE Trans ASSP*, vol.38, July 1990, 1257-1265.

[31] J.K.TUGNAIT, "Approaches to FIR System Identification with Noisy Data using Higher-Order Statistics", *Proc Workshop Higher-Order Spectral Analysis*, June 1989, Vail, 13-18.

[32] J.K.TUGNAIT, "Identification of Linear Stochastic Systems via Second and Fourth Order Cumulant Matching", *IEEE Trans Inf Theory*, vol.33, May 1987, 393-407.

[33] D.L.WALLACE, "Asymptotic Approximations to Distributions", *Annals Math Statistics*, 1958, 29, 635-654.

Appendices

◆Proof of lemma (22)

Mutual independence of the components of \mathbf{x} yields

$$(A-1) \quad \varphi_{\mathbf{x}}(\mathbf{u}) = \prod_i \varphi_{x_i}(u_i).$$

On the other hand, the relation $\mathbf{z} = \mathbf{B} \mathbf{x}$ implies

$$(A-2) \quad \varphi_{\mathbf{z}}(\mathbf{u}) = \varphi_{\mathbf{x}}(\mathbf{B}^* \mathbf{u}).$$

A characteristic function is always continuous. Therefore, since it takes the value 1 at the origin, there exist a neighborhood of zero, U , in which $\varphi_{\mathbf{z}}(\mathbf{u})$ and $\varphi_{\mathbf{x}}(\mathbf{u})$ are strictly positive. For convenience, denote by χ the logarithm of φ for $\mathbf{u} \in U$. Then (A-1) and (A-2) give together

$$(A-3) \quad \chi_{\mathbf{z}}(\mathbf{u}) = \sum_i \chi_{x_i}(\mathbf{B}_{:i}^* \mathbf{u}),$$

where $\mathbf{B}_{:i}$ denotes the i th column of \mathbf{B} . Now suppose there exist two non-zero entries B_{pi} and B_{qi} in \mathbf{B} . Then from (A-3), and by pairwise independence of z_p and z_q , we may write:

$$\sum_i \chi_{x_i}(\mathbf{B}_{pi} u_p) + \chi_{x_i}(\mathbf{B}_{qi} u_q) = \sum_i \chi_{x_i}(\mathbf{B}_{pi} u_p + \mathbf{B}_{qi} u_q).$$

Now from a result of Darmois [14], we know that all χ_{x_i} for which both B_{pi} and B_{qi} are non-zero are polynomials of degree at most 2. \diamond

◆Proof of theorem (23)

Implications (iii) \Rightarrow (ii) and (ii) \Rightarrow (i) are quite obvious. We shall prove the last one, namely (i) \Rightarrow (iii). Assume \mathbf{z} has pairwise independent components, and suppose \mathbf{C} is not of the form $\Lambda \mathbf{P}$. Since \mathbf{C} is orthogonal, it has necessarily two nonzero entries in at least two different columns. Then by applying the lemma twice, \mathbf{x} has at least two gaussian components, which is contrary to our hypothesis. \diamond

◆Proof of lemma (41)

The matrix ${}^2\mathbf{Q}$ is bistochastic (i.e. the sum of its entries in any row or column is equal to 1). But from Birkhoff theorem, the set of bistochastic matrices is a convex polyhedron whose vertices are permutations. Thus ${}^2\mathbf{Q}$ can be decomposed as

$$(A-5) \quad {}^2\mathbf{Q} = \sum_s \alpha_s \mathbf{P}_s, \quad \alpha_s \geq 0, \quad \sum_s \alpha_s = 1.$$

Then we have the inequality

$$(A-6) \quad \|{}^2\mathbf{Q} \mathbf{u}\| \leq \sum_s \alpha_s \|\mathbf{P}_s \mathbf{u}\| = \|\mathbf{u}\|. \diamond$$

◆Proof of theorem (42)

Since \mathbf{Q} is unitary, we have $|Q_{ij}| \leq 1$, and consequently $|{}^1Q_{ij}| \leq |{}^2Q_{ij}|$. And applying the triangular inequality

$$(A-7) \quad \sum_{i,j,k} Q_{ki}^* Q_{kj}^* u_i u_j \leq \sum_{i,j,k} Q_{ki}^2 Q_{kj}^2 |u_i| |u_j|.$$

Then using the lemma we get

$$(A-8) \quad \|{}^1\mathbf{Q} \mathbf{u}\|^2 \leq \|{}^2\mathbf{Q} \mathbf{u}\|^2 \leq \|\mathbf{u}\|^2. \diamond$$

The details of the proofs are deferred to a full paper [7].

Acknowledgement

The author wishes to thank A.Groenenboom and S.Proserpi for their proofreading of the paper. The reference [25] indicated by Dr Swami some years ago was also helpful.